

PENERAPAN PARTIAL LEAST SQUARES PADA DATA GINGEROL

Margaretha Ohyver

Jurusan Matematika dan Statistik, Fakultas Sains dan Teknologi, Universitas Bina Nusantara
Jln. K.H. Syahdan No. 9, Palmerah, Jakarta Barat 11480
ethaohyver@binus.ac.id

ABSTRACT

Multivariate calibration model aims to predict the expensive measures obtained by using the measures of a cheap and easy. There are several problems that often occur in the model calibration, among others, and multikolinear. To overcome these problems we used partial least squares method (PLS). The study was conducted to apply the PLS method on the data gingerol. Based on research conducted with the two components of the model obtained with the diversity of variable Y at 83.8032% and the diversity of variable X equal to 100%, and obtained for $R^2 = 83.8\%$ and $RMSE = 0.100891$ calibration data group and $R^2 = 84.2\%$ and $RMSEP = 0.199939$ for the validation data.

Keywords: gingerol, multivariate calibration, partial least squares

ABSTRAK

Model multivariate calibration bertujuan untuk menduga ukuran-ukuran yang mahal diperoleh dengan menggunakan ukuran-ukuran yang murah dan mudah. Ada beberapa masalah yang sering terjadi pada pemodelan kalibrasi, diantaranya ($n < p$) dan multikolinear. Untuk mengatasi permasalahan tersebut maka digunakan metode partial least squares (PLS). Penelitian dilakukan untuk menerapkan metode PLS pada data gingerol. Berdasarkan penelitian yang dilakukan diperoleh model dengan 2 komponen dengan keragaman peubah Y sebesar 83,8032% dan keragaman peubah X sebesar 100% serta diperoleh untuk $R^2 = 83,8\%$ dan $RMSE = 0,100891$ kelompok data kalibrasi dan $R^2 = 84,2\%$ dan $RMSEP = 0,199939$ untuk kelompok data validasi.

Kata kunci: gingerol, multivariate calibration, partial least squares.

PENDAHULUAN

Model *multivariate calibration* merupakan bagian dari *chemometric*. Di mana *chemometric* merupakan disiplin kimia yang menggabungkan metode matematika dan statistika dengan kimia. *Multivariate calibration* bertujuan untuk menemukan model yang dapat digunakan untuk menduga ukuran-ukuran yang mahal diperoleh dengan menggunakan ukuran-ukuran yang murah dan mudah diperoleh secara tepat dan akurat. Secara umum *multivariate calibration* menggunakan formula matematika untuk menduga informasi pada Y, yaitu ukuran yang mahal, yang tidak diketahui berdasarkan informasi pada X, yaitu ukuran yang murah, yang tersedia. Formula matematika yang disebut model pada prinsipnya dibagi menjadi dua komponen, yaitu komponen struktur dan komponen sisaan. Komponen struktur adalah komponen yang menggambarkan variasi sistematis, sedangkan komponen sisaan adalah komponen yang menggambarkan perbedaan antara data dan komponen struktur (Martens dan Naes, 1989). Terdapat beberapa masalah dalam model *multivariate calibration*, diantaranya banyaknya pengamatan lebih kecil daripada banyaknya peubah ($n < p$) dan multikolinear. Salah satu metode yang dapat digunakan untuk mengatasi kedua masalah di atas adalah metode *partial least squares*.

Partial Least Squares (PLS) merupakan perpaduan antara Analisis Komponen Utama (AKU) dan regresi linear ganda (Abdi, 2003). PLS membentuk peubah bebas yang baru yang disebut faktor, peubah laten, atau komponen, di mana masing-masing komponen yang terbentuk merupakan kombinasi linear dari peubah-peubah bebas. Metode PLS mempunyai persamaan dengan *Principal Component Regression* (PCR). Persamaannya adalah keduanya menggunakan komponen sebagai peubah bebas. Adapun perbedaannya adalah komponen pada PCR hanya ditentukan dari peubah bebas, sedangkan komponen untuk PLS ditentukan oleh peubah bebas dan peubah respon. Tujuan utama dari PLS adalah membentuk komponen yang dapat menangkap informasi dari peubah bebas untuk menduga peubah respon (Hoskuldsson dalam Garthwaite, 1994).

Jika ($n < p$), maka metode metode kuadrat terkecil tidak dapat digunakan (Naes, Isaksson, dkk., 2002). Hal ini dikarenakan matriks $\mathbf{X}^T\mathbf{X}$ singular. Sebaliknya, PLS dapat digunakan untuk kasus ($n < p$). Regresi PLS didasarkan pada dekomposisi komponen:

$$\mathbf{Y} = \mathbf{TQ}' + \mathbf{F} \quad (1)$$

$$\mathbf{X} = \mathbf{TP}' + \mathbf{E} \quad (2)$$

dengan \mathbf{T} adalah matriks komponen, \mathbf{P} dan \mathbf{Q} adalah matriks *loading* X dan Y, \mathbf{E} dan \mathbf{F} adalah vektor *error* (Boulesteix dan Strimmer, 2006).

Metode PLS dapat dipandang sebagai metode yang membentuk matriks komponen \mathbf{T} sebagai transformasi linear dari \mathbf{X} .

$$\mathbf{T} = \mathbf{XW} \quad (3)$$

dengan \mathbf{W} adalah matriks *weights* (bobot). Persamaan (3) dapat dituliskan sebagai berikut.

$$T_1 = w_{11}x_1 + w_{21}x_2 + \dots + w_{p1}x_p$$

$$T_2 = w_{12}x_1 + w_{22}x_2 + \dots + w_{p2}x_p$$

$$\vdots$$

$$T_c = w_{1c}x_1 + w_{2c}x_2 + \dots + w_{pc}x_p$$

Komponen-komponen kemudian digunakan untuk pendugaan, dengan menggantikan \mathbf{X} sehingga diperoleh penduga kuadrat terkecil:

$$\hat{\mathbf{Q}}' = (\mathbf{T}'\mathbf{T})^{-1}\mathbf{T}'\mathbf{Y} \quad (4)$$

Metode PLS diawali dengan mentransformasikan peubah bebas (X) dan peubah respon (Y).

$$x_{ij}^* = \frac{x_{ij} - \bar{x}_j}{S_{x_j}} \quad (5)$$

$$y_i^* = \frac{y_i - \bar{y}}{S_y} \quad (6)$$

di mana \bar{x}_j adalah rata-rata nilai x_{ij} , \bar{y} adalah rata-rata nilai y . Sedangkan S_{x_j} dan S_y adalah simpangan baku x_j dan y , yaitu:

$$S_{x_j} = \sqrt{\frac{\sum (x_{ij} - \bar{x}_j)^2}{n-1}}, \quad (7)$$

$$S_y = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n-1}}. \quad (8)$$

Algoritma PLS

Algoritma PLS terdiri dari sebagai berikut. Pertama, transformasi peubah X dan Y menjadi X^* dan Y^* . Kedua, mengambil nilai awal vektor $\mathbf{u} = \mathbf{Y}^*$. Ketiga, menentukan bobot \mathbf{X}^* , dengan persamaan $\mathbf{w}^T = \frac{\mathbf{u}^T \mathbf{X}^*}{\mathbf{u}^T \mathbf{u}}$. Keempat, menentukan $\mathbf{w}^* = \frac{\mathbf{w}_{lama}^T}{\|\mathbf{w}_{lama}^T\|}$. Kelima, menentukan faktor skor \mathbf{X}^* ,

dengan persamaan $\mathbf{t} = \frac{\mathbf{X}^* \mathbf{w}^*}{\mathbf{w}^{*T} \mathbf{w}^*}$. Keenam, menentukan bobot \mathbf{Y}^* , dengan persamaan $\mathbf{c}^T = \frac{\mathbf{t}^T \mathbf{Y}^*}{\mathbf{t}^T \mathbf{t}}$.

Ketujuh, menentukan $\mathbf{c}^* = \frac{\mathbf{c}_{lama}^T}{\|\mathbf{c}_{lama}^T\|}$. Kedelapan, menentukan skor \mathbf{Y}^* , dengan persamaan $\mathbf{v} = \frac{\mathbf{Y}^* \mathbf{c}^*}{\mathbf{c}^{*T} \mathbf{c}^*}$.

Kesembilan, menentukan $\mathbf{b} = \frac{\mathbf{v}^T \mathbf{t}}{\mathbf{t}^T \mathbf{t}}$. Kesepuluh, menentukan faktor *loading* untuk \mathbf{X} , dengan

persamaan $\mathbf{p}^T = \frac{\mathbf{t}^T \mathbf{X}^*}{\mathbf{t}^T \mathbf{t}}$. Kesebelas, menentukan $\mathbf{E}_h = \mathbf{E}_{h-1} - \mathbf{t}_h \mathbf{p}_h^T$; $\mathbf{E}_0 = \mathbf{X}^*$.

Keduabelas, menentukan $\mathbf{F}_h = \mathbf{F}_{h-1} - \mathbf{b}_h \mathbf{t}_h \mathbf{c}_h^T$; $\mathbf{F}_h = \mathbf{Y}^*$.

Untuk memeriksa kebaikan modelnya, digunakan statistik *Prediction Sum of Squares* (PRESS). Persamaannya adalah:

$$\text{PRESS} = \sum_{i=1}^n (y_i - \hat{y}_{i,-i})^2 \quad (9)$$

Dengan y_i adalah nilai peubah respon pada pengamatan ke- i , dan $\hat{y}_{i,-i}$ adalah nilai dugaan y_i tanpa pengamatan ke- i .

Model dengan nilai PRESS terkecil mengindikasikan kecilnya galat pendugaan dalam model sehingga suatu model dikatakan lebih baik jika nilai PRESS yang dihasilkan relatif lebih kecil.

Prosedur PRESS ditempuh melalui cara menyisihkan satu pengamatan, menduga modelnya dari amatan yang ada lalu menduga pengamatan yang disisihkan sebelumnya serta menghitung kuadrat selisih antara pengamatan dan dugaan. Prosedur ini dilakukan untuk setiap pengamatan.

Langkah terakhir di dalam proses pembentukan model adalah validasi model regresi yang terpilih. Terdapat beberapa metode validasi, diantaranya membagi data menjadi dua bagian. Data bagian pertama, dinamakan *model building set*, digunakan untuk membangun model. Yang kedua, dinamakan *validation or prediction set*, digunakan untuk menguji model (Neter, dkk., 1990).

Salah satu ukuran yang dapat digunakan untuk validasi model adalah dengan menghitung *root mean squared error prediction* (RMSEP):

$$RMSEP = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2} \quad (10)$$

dengan Y_i adalah nilai peubah respon dalam kasus validasi ke- i , \hat{Y}_i adalah nilai dugaan dalam kasus validasi ke- i , dan n adalah banyaknya pengamatan dalam kasus validasi. Model dikatakan valid jika memiliki nilai RMSEP yang kecil.

Pendugaan model kalibrasi dengan menggunakan metode PLS telah banyak dilakukan. Di antaranya aplikasi PLS dalam penentuan *Chlorogenic Acid* pada sampel tanaman (Shao dan Zhuang, 2004), dan aplikasi PLS yang didasarkan pada resolusi spektra *Near Infrared* (NIR) yang berbeda (Chung, Choi, Choo, dan Lee, 2004). Untuk penelitian kali ini, data yang akan digunakan adalah data *gingerol* pada rimpang jahe.

Jahe merupakan salah satu dari beberapa tanaman yang digunakan secara tradisional sebagai obat rematik, demam, radang, dan lain-lain. Rimpang jahe mengandung dua bagian utama, yaitu *volatil* (minyak esensial) yang memberikan aroma dan *gingerol* yang merupakan pembawa rasa pedas. Kandungan *gingerol* yang cukup tinggi pada rimpang jahe, menyebabkan jahe memiliki peranan yang sangat penting. Peranan penting yang dimaksud adalah peranan dalam dunia pengobatan baik pengobatan tradisional atau skala industri dengan memanfaatkan kemajuan teknologi.

Data *gingerol* ini sebelumnya telah digunakan pada beberapa penelitian. Diantaranya digunakan untuk penerapan metode *neural network* dan *principal component* (Atok dan Notodiputro, 2004), untuk penerapan PCR (Arnita, 2005), dan untuk penerapan Transformasi Wavelet Diskrit (TWD) dengan menggunakan *mother wavelet* Haar dan PCR (Sunaryo, 2005). Pada data *gingerol*, hasil yang diperoleh dengan menggunakan TWD-RKU lebih baik daripada yang diperoleh tanpa ditransformasi wavelet. Hal ini ditinjau dari kriteria *Root Mean Squared Error Prediction* (RMSEP). RMSEP untuk TWD-RKU adalah 0,1072, sedangkan RMSEP untuk RKU adalah 0,1430.

Data *gingerol* merupakan salah satu contoh data yang mempunyai pengamatan lebih kecil daripada banyaknya peubah ($n < p$) dan multikolinear. Karena tidak tertutup kemungkinan akan ditemukan data sejenis ini, maka perlu metode-metode untuk mengatasinya sehingga tujuan dari penelitian ini adalah menerapkan metode PLS untuk mengatasi masalah ($n < p$) dan multikolinear.

METODE PENELITIAN

Data yang digunakan dalam penelitian ini merupakan data sekunder yang diperoleh dari peneliti terdahulu (Sunaryo, 2005). Data ini berupa data pengamatan senyawa aktif *gingerol* pada rimpang jahe. Alat yang digunakan untuk memperoleh kandungan senyawa aktif *gingerol* adalah HPLC dan FTIR. Hasil pengukuran bilangan gelombang dengan FTIR berupa data spektra %

transmitan. Setiap bentuk spektrum % transmitan dari FTIR akan mencerminkan gugus fungsi yang terdapat pada senyawa (dalam hal ini *gingerol*) dari suatu sampel rimpang jahe. Dengan menggunakan FTIR dihasilkan data spektra % transmitan sebanyak 1866 titik pada bilangan gelombang 4000-200 cm^{-1} yang mencerminkan kadar *gingerol*. % transmitan sebanyak 1866 titik ini digunakan sebagai peubah bebas (X), sedangkan kadar senyawa aktif hasil pengukuran dari HPLC sebagai peubah respon.

Dari penelitian sebelumnya, diketahui bahwa rimpang jahe mengalami masa simpan yang berbeda-beda pada tiap sampel (Sunaryo, 2005; Arnita, 2005). Lama masa simpan ini ternyata berpengaruh terhadap kadar *gingerol* yang dihasilkan. Oleh karena itu peubah *dummy* akan diikutkan dalam pendugaan model yang mencerminkan kelompok lama masa simpan.

Lama masa simpan dikategorikan menjadi dua, yaitu masa simpan lama (kode 1) dan masa simpan sebentar (kode 0). Yang dimaksud dengan masa simpan lama adalah masa simpan yang lebih dari 3 bulan dan masa simpan sebentar adalah masa simpan yang kurang dari 3 bulan.

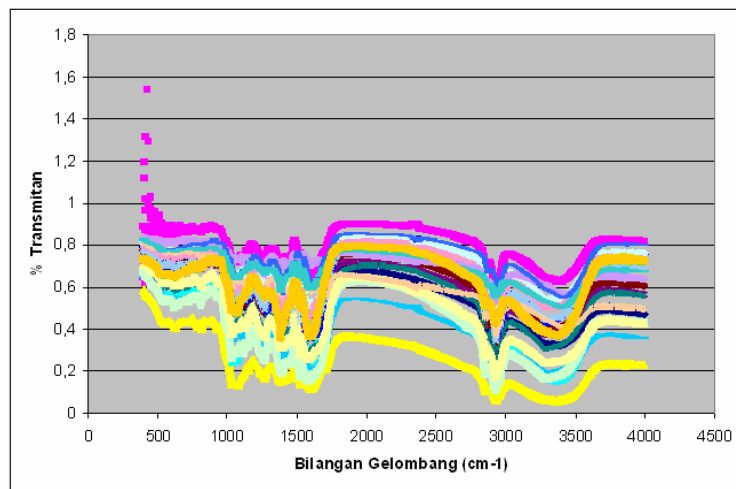
Tabel 1 Dua Puluh Sampel Rimpang Jahe

No	Sampel	Kadar <i>Gingerol</i>	Masa Simpan	Kode
1	Suharsono1	0,63	10 bulan	1
2	Suharsono2	0,53	10 bulan	1
3	Dukuh1	0,72	10 bulan	1
4	Dukuh2	0,78	10 bulan	1
5	Suparno1	0,58	10 bulan	1
6	Suparno2	0,53	10 bulan	1
7	Karyo1	0,52	10 bulan	1
8	Karyo2	0,54	10 bulan	1
9	Haryono1	0,79	10 bulan	1
10	Haryono2	0,78	10 bulan	1
11	Mulyono1	0,63	10 bulan	1
12	Mulyono2	0,63	10 bulan	1
13	Sugandi1	0,78	10 bulan	1
14	Sugandi2	0,79	10 bulan	1
15	Majalengka1	1,26	< 3 bulan	0
16	Majalengka2	1,6	< 3 bulan	0
17	Balitro1	1,18	< 3 bulan	0
18	Balitro2	1,14	< 3 bulan	0
19	Bogor1	1,24	< 3 bulan	0
20	Bogor2	1,07	< 3 bulan	0

Penelitian ini dilakukan dengan langkah-langkah sebagai berikut. Pertama, membagi data menjadi dua bagian. Lima pengamatan yang dipilih secara acak digunakan untuk validasi model. Lima belas pengamatan digunakan untuk membentuk model *multivariate calibration*. Kedua, membentuk model *multivariate calibration* dengan menggunakan metode PLS. Peubah *dummy* diikutkan. Ketiga, menentukan banyaknya komponen pada metode PLS dengan menggunakan akar rata-rata PRESS. Keempat, menghitung RMSE dan R^2 untuk kelompok data kalibrasi dan kelompok data validasi.

HASIL DAN PEMBAHASAN

Dari dua puluh pengamatan, akan digunakan lima belas pengamatan untuk pengembangan model kalibrasi dan lima pengamatan untuk validasi model. Peubah bebas yang akan digunakan adalah 1867, karena peubah *dummy*, yaitu lama masa simpan rimpang jahe diikutsertakan. Karena $n < p$, maka koefisien-koefisien dugaan yang dihasilkan untuk masing-masing parameter akan beragam.



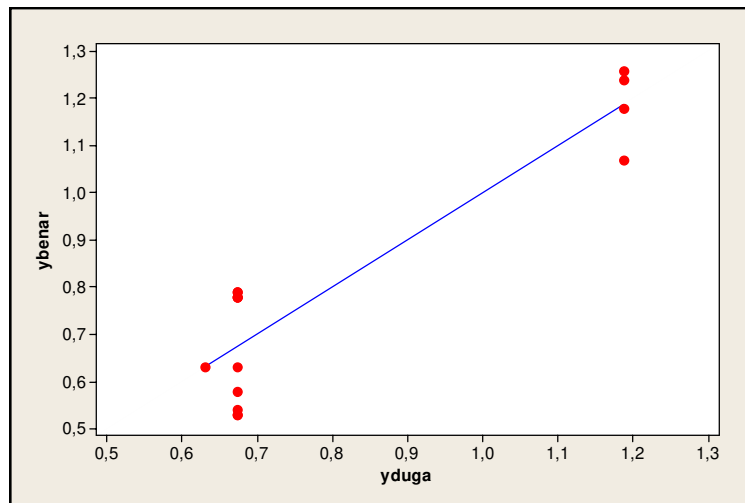
Gambar 1 Spektra % Transmittan 1866 Titik

Hasil metode PLS beserta nilai keragaman % transmittan (X) dan hasil pengukuran konsentrasi senyawa aktif *gingerol* rimpang jahe menggunakan HPLC (Y) dapat dilihat pada Tabel 2. Berdasarkan tabel tersebut, terlihat bahwa kenaikan tertinggi (80,1459 %) untuk nilai keragaman peubah Y diperoleh untuk 2 komponen. Sedangkan untuk peubah X , persentase keragamannya telah mencapai 100% untuk 2 komponen.

Tabel 2 Banyaknya Komponen Hasil Metode PLS

Banyaknya Komponen	X		Y		ARP
	%keragaman	%kumulatif	%keragaman	%kumulatif	
1	99,9472	99,9472	3,6574	3,6574	1491,05
2	0,0528	100,0000	80,1459	83,8032	419,46
3	0,0000	100,0000	0,9051	84,7083	4587,61
4	0,0000	100,0000	4,6441	89,3525	6064,04
5	0,0000	100,0000	1,9416	91,2940	5903,39
6	0,0000	100,0000	0,7553	92,0493	5151,60
7	0,0000	100,0000	1,5210	93,5703	5545,45
8	0,0000	100,0000	1,5600	95,1303	5607,66
9	0,0000	100,0000	3,3492	98,4795	6383,45
10	0,0000	100,0000	0,7651	99,2446	7949,03
11	0,0000	100,0000	0,3866	99,6312	8799,50
12	0,0000	100,0000	0,2667	99,8979	8799,50
13	0,0000	100,0000	0,1018	99,9997	8799,50
14	0,0000	100,0000	0,0003	100,0000	8799,50
15	0,0000	100,0000	0,0000	100,0000	8799,50

Nilai akar rata-rata PRESS (ARP) yang dihasilkan adalah sangat tinggi. Untuk nilai ARP yang tertinggi adalah 8799,50 yang diperoleh untuk 11, 12, 13, 14, dan 15 komponen. Sedangkan nilai ARP terkecil adalah 419,46 yang diperoleh untuk 2 komponen. Jika didasarkan pada nilai ARP terkecil, maka yang digunakan adalah sebanyak 2 komponen. Dengan 2 komponen, model telah dapat menjelaskan sekitar 83,8032 % dari keragaman peubah Y dan 100 % dari keragaman peubah X .



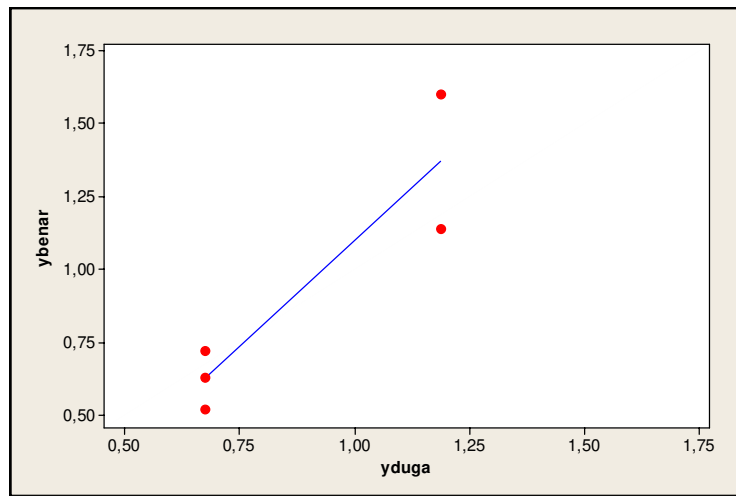
Gambar 2 Plot Y dengan \hat{Y} untuk Kelompok Data Kalibrasi dengan Metode PLS

Berdasarkan komponen yang telah diperoleh, maka akan dicari nilai-nilai dugaan untuk peubah Y . Nilai dugaan peubah Y diperoleh dengan mengalikan koefisien-koefisien dugaan dengan nilai-nilai peubah X . Nilai-nilai dugaan peubah Y untuk kelompok data kalibrasi dan kelompok data validasi dapat dilihat pada Tabel 3. Plot untuk Plot Y dengan \hat{Y} untuk kelompok data kalibrasi dan kelompok data validasi dapat dilihat pada Gambar 2 dan 3.

Tabel 3 Nilai Y dan \hat{Y} untuk *Gingerol Rimpang Jahe* dengan Metode PLS

Kelompok Data Kalibrasi		Kelompok Data Validasi	
Kadar Gingerol dari HPLC (%)	Dugaan (%)	Kadar Gingerol dari HPLC (%)	Dugaan (%)
0,63	0,62998	0,72	0,67301
0,53	0,67301	0,52	0,67300
0,78	0,67300	0,63	0,67300
0,58	0,67300	1,60	1,18750
0,53	0,67300	1,14	1,18750
0,54	0,67300		
0,79	0,67300		
0,78	0,67300		
0,63	0,67300		
0,78	0,67300		
0,79	0,67300		
1,26	1,18750		
1,18	1,18750		
1,24	1,18750		
1,07	1,18750		

Berdasarkan nilai-nilai dugaan tersebut diperoleh untuk kelompok data kalibrasi, $R^2 = 83,8 \%$ dan $RMSE = 0,100891$. Sedangkan untuk kelompok data validasi diperoleh $R^2 = 84,2 \%$ dan $RMSEP = 0,199939$.



Gambar 3 Plot Y dengan \hat{Y} untuk Kelompok Data Validasi dengan Metode PLS

SIMPULAN

Dari penelitian yang dilakukan maka diperoleh kesimpulan sebagai berikut. Pertama, penerapan metode PLS pada data *gingerol* diperoleh model dengan 2 komponen dengan keragaman peubah Y sebesar 83,8032 % dan keragaman peubah X sebesar 100%. Kedua, dengan 2 komponen diperoleh untuk $R^2 = 83,8 \%$ dan $RMSE = 0,100891$ kelompok data kalibrasi dan $R^2 = 84,2 \%$ dan $RMSEP = 0,199939$ untuk kelompok data validasi.

Jika dilihat dari kriteria R^2 dan $RMSE$, baik untuk data kalibrasi maupun data validasi, maka model dengan 2 komponen bisa dikatakan baik. Tetapi jika dibandingkan dengan penelitian yang dilakukan oleh Sunaryo, maka R^2 dan $RMSE$ dengan metode TWD-PCR masih lebih baik. Sunaryo melakukan penelitiannya dengan mentransformasi data *gingerol* menggunakan wavelet. Hal ini dilakukan dengan pertimbangan adanya *noise* sehingga untuk penelitian selanjutnya bisa dilakukan dengan menggunakan gabungan antara metode wavelet dan PLS.

DAFTAR PUSTAKA

- Abdi, H. (2003). Partial least squares regression. *Encyclopedia of Social Sciences Research Methods* (online), 1-7. Retrieved from <http://www.utdallas.edu/~herve>.
- Arnita. (2005). *Koreksi pencaran pada data kalibrasi rimpang jahe (Zingiber officinale)*. Tesis tidak diterbitkan, Bogor: Program Pascasarjana, Institut Pertanian Bogor.
- Atok, R. M., dan Notodiputro, K. A. (2004). Metode NN (Neural Network) dengan principle component sebagai pre-processing pada data. *Proceeding Seminar Nasional Statistika*, Bogor: Institut Pertanian Bogor.
- Boulesteix, A., and Strimmer, K. (2006). *Partial least squares: A versatile tool for the analysis of high-dimensional genomic data* (online). Retrieved from

<http://www.slcmsr.net/boulesteix/papers/review>.

- Chung, H. *et al.* (2004). Investigation of partial least squares calibration performance based on different resolutions of near infrared spectra. *Bull. Korean Chem. Soc*, 25 (5), 647-651.
- Garthwaite, P. H. (1994). An interpretation of partial least squares. *Journal of the American Statistical Association*, 89, 122-127.
- Martens, H., and Naes, T. (1989). *Multivariate calibration*, New York: John Wiley & Sons, Inc.
- Naes, T. *et al.* (2002). *Multivariate calibration and classification*, Chichester: NIR Publications.
- Neter, J., Wasserman, W., and Kutner, M. H. (1990). *Applied linear statistical models*, Illinois: Irwin.
- Shao, X., and Zhuang, Y. (2004). Determination of chlorogenic acid in plant samples by using near-infrared spectrum with wavelet transform preprocessing. *Analytical Sciences*, 20, 451-454.
- Sunaryo, S. (2005). *Model kalibrasi dengan transformasi wavelet sebagai metode pra-pemrosesan*. Disertasi tidak diterbitkan. Bogor: Program Pascasarjana, Institut Pertanian Bogor.